



Heriot-Watt University
Research Gateway

Language and domain aware lightweight ontology matching

Citation for published version:

Bella, G, Giunchiglia, F & McNeill, F 2017, 'Language and domain aware lightweight ontology matching', *Journal of Web Semantics*, vol. 43, pp. 1-17. <https://doi.org/10.1016/j.websem.2017.03.003>

Digital Object Identifier (DOI):

[10.1016/j.websem.2017.03.003](https://doi.org/10.1016/j.websem.2017.03.003)

Link:

[Link to publication record in Heriot-Watt Research Portal](#)

Document Version:

Peer reviewed version

Published In:

Journal of Web Semantics

General rights

Copyright for the publications made accessible via Heriot-Watt Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

Heriot-Watt University has made every reasonable effort to ensure that the content in Heriot-Watt Research Portal complies with UK legislation. If you believe that the public display of this file breaches copyright please contact open.access@hw.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Contents lists available at ScienceDirect

Web Semantics: Science, Services and Agents on the World Wide Web

journal homepage: www.elsevier.com/locate/websem

Language and domain aware lightweight ontology matching

Gábor Bella^{a,*}, Fausto Giunchiglia^a, Fiona McNeill^b^a University of Trento, via Sommarive 5, 38123 Trento, Italy^b Heriot-Watt University, Edinburgh EH14 4AS, Scotland, United Kingdom

ARTICLE INFO

Article history:

Received 24 August 2016

Received in revised form 15 March 2017

Accepted 26 March 2017

Available online xxxx

Keywords:

Cross-lingual matching

Multilingual matching

Domains

Ontology matching

Semantic matching

Machine translation

ABSTRACT

Concepts and relations in ontologies and in other knowledge organisation systems are usually annotated with natural language labels. Most ontology matchers rely on such labels in element-level matching techniques. State-of-the-art approaches, however, tend to make implicit assumptions about the language used in labels (usually English) and are either domain-agnostic or are built for a specific domain. When faced with labels in different languages, most approaches resort to general-purpose machine translation services to reduce the problem to monolingual English-only matching. We investigate a thoroughly different and highly extensible solution based on *semantic matching* where labels are parsed by multilingual natural language processing and then matched using language-independent and domain aware background knowledge acting as an interlingua. The method is implemented in NuSM, the language and domain aware evolution of the SMATCH semantic matcher, and is evaluated against a translation-based approach. We also design and evaluate a fusion matcher that combines the outputs of the two techniques in order to boost precision or recall beyond the results produced by either technique alone.

© 2017 Elsevier B.V. All rights reserved.

1. Introduction

Ontologies and other knowledge organisation systems, while usually serving a purpose of standardisation or generalisation, stem from local needs and practices. By *local* we understand *within an administrative unit* such as a country or a region as well as *within an application domain* such as medicine or transport. Accordingly, ontologies tend to target specific domains and the labels annotating their elements – concepts, relations, metadata – tend to be expressed in the local language. This is especially true for *lightweight ontologies* [1]: classification hierarchies, taxonomies, and other tree-structured data schemas widely used around the world as simple, well-understood, semi-formal resources for knowledge organisation. Such resources often play normative roles on the national level in public services, industry, or commerce, as a means for classification (of documents, books, open data, commercial products, web pages, etc.) as well as being sources of shared vocabularies for actors cooperating in a given domain.

Ontology matching [2] is a process that creates and maintains alignments between elements of two ontologies covering overlapping areas of knowledge. We define *language aware* or *multilingual matching* as a type of ontology matching where a multilingual setting is explicitly assumed, i.e., the matcher is capable of dealing with ontologies expressed in multiple languages. Likewise, we

define *domain aware* matching as capable of dealing with domain-specific knowledge and domain terms with specialised meanings.

Activities on supra-national levels such as international trade and mobility need to rely on the interoperability and integration of knowledge organisation resources across countries, languages, and sometimes across domains. *Cross-lingual matching* is a specific case of language aware matching when ontologies in different languages need to be aligned. Likewise, *cross-domain matching* is used to match ontologies pertaining to different domains of knowledge. An example of a simultaneously cross-lingual and cross-domain matching problem is the case of *cross-border emergency response* where responders from different countries and from different domains (geography, geology, medicine, police, military, transportation, etc.) need to share data. In [3] we apply the domain aware matching approach presented in this paper to this particular use case.

State-of-the-art cross-lingual matchers invariably use translation-based techniques – most often online machine translation services from Microsoft or Google – in order to reduce the problem of multilingualism to the well-researched problem of monolingual English-to-English matching [2,4–8, p. 105]. With the constant improvement of such services, translation-based matchers are able to provide usable results and are able to deal with a wide range of languages. State-of-the-art machine translators today mainly use statistical methods and are trained on large amounts of *bilingual parallel* or *comparable corpora* for each language pair they support.

* Corresponding author.

E-mail address: gabor.bella@unitn.it (G. Bella).

A known problem of statistical machine translation, however, is the decrease of translation accuracy on corpora significantly different from those on which the system was trained. This typically happens on domain classifications and ontologies that contain specialised terminology. The adaptation of a statistical system to a new domain requires re-training on corpora extended with a significant amount of domain-specific text (ideally bilingual parallel corpora that are hard to find). At the same time, the systems typically used by ontology matchers are on-line commercial services (such as Bing and Google Translate) that, while offering the best available translation quality, are not adaptable or customisable by the user.

The shortness of labels typically found in ontologies is another difficulty that state-of-the-art approaches face, as the sparseness of textual context within labels makes the translation task more error-prone. Furthermore, the often non-standard orthography and syntax of ontology labels – that we described in [9] as a form of specialised *block language* – makes label parsing even harder.

In this paper we introduce a different approach to language and domain aware matching, so far hardly investigated, that does not rely on external translation tools. The method is based on combining two types of resources: on the one hand, multilingual natural language processing tools that are adapted to the block language of structured data and, on the other hand, off-line *multilingual lexical databases* connecting words and expressions of natural language to language-independent but domain-aware meanings.

The work described in this paper is motivated by the following considerations. Firstly, while both approaches evoked above are resource-intensive, the types of resources they feed on are markedly different: on the one hand, machine translation requires large amounts of bilingual parallel or comparable corpora relevant to the target domain, on the other hand, our approach uses lexical, terminological, and NLP resources for each supported language. In both cases, a wide range of resources are already available on the web. Based on their availability and conditions of use, for specific use cases one approach or the other may be more cost-effective or faster to implement. Our knowledge-based label matching approach can thus be seen as an alternative when no good-quality language or domain-specific machine translator is available. Secondly, we are interested in comparing the strengths and weaknesses of the two approaches, which turn out to be rather complementary. Our evaluations use two machine translation systems: Google Translate, currently the best available online translator, and Apertium, which is free and can also be used off-line. We conduct evaluations on three language pairs: English–Spanish, English–Italian, and Spanish–Italian. Finally, based on the complementarity of the two approaches we investigate the idea of combining them – using multilingual lexical resources on the one hand and machine translation on the other hand – into a single matcher. The resulting system, as demonstrated by our evaluation results, clearly outperforms either method alone.

The result of our work is implemented in *NuSMATCH* (NuSM for short), an upcoming release of the open-source SMATCH system [10] with built-in capabilities for language and domain aware matching.

The rest of the paper is organised as follows. After a brief reminder of semantic matching in Section 2, Section 3 introduces language and domain aware matching. Section 4 presents the core multilingual and multidomain resource serving as background knowledge for our matcher. Sections 5 and 6 deal with interpreting labels, explaining how NuSM performs language aware parsing and domain aware sense disambiguation, respectively. Section 7 describes our cross-lingual label matching techniques. Section 8 provides mechanisms for extending NuSM by languages and domain terms. Section 9 proposes a fusion matcher that combines NuSM with translation-based matching. Section 10 provides evaluations for NuSM both in comparison to and in combination with translation-based techniques. Section 11 presents related work and, finally, Section 12 draws the overall conclusions.

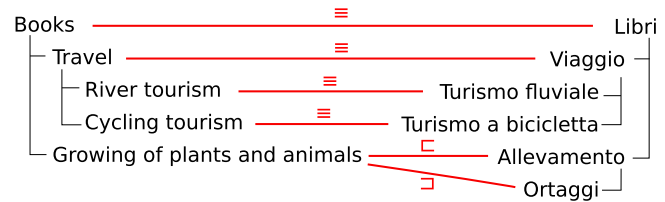


Fig. 1. Example English and Italian classifications of books, with some example mapping relations.

2. Semantic matching

In this section we provide an overview of the notions of *lightweight ontology*, *semantic matching*, and of the principal aspects of language and domain awareness in NuSM.

2.1. Principles of semantic matching

NuSM is designed as a multilingual and domain aware extension of the SMATCH (English-only) semantic matcher [10,11]. SMATCH was specifically designed for matching semi-formal knowledge organisation schemes such as classification hierarchies (as opposed to formal ontologies expressed, e.g., in OWL) that we believe are the main subject of most real-world multilingual matching applications. As shown in [1], such classifications typically have the following properties:

- a tree structure;
- nodes are expressed as short and well-formed natural language labels;
- classification semantics are implied for nodes and edges.

As an example of classification semantics, in a classification of newspaper articles the extension of a node *Literature* are articles about literature, while a node *Italy* under *Literature* classifies articles on Italian literature.

Matching is qualified as *semantic* for three reasons:

- it is performed using logical inference on propositional description logic formulas that capture the meanings of natural-language classification labels in a formal way;
- atoms of the formulas are concepts and are not based on the surface forms of words;
- the ontology node mappings output by the matcher are ‘meaningful’ description logic relations of equivalence, subsumption, and disjointness as opposed to similarity scores.

These operators are important in several real-world matching tasks such as data integration or query answering. For example, in a schema matching operation between query and database attributes, in response to a query on a ‘*phone_number*’ a database may return a ‘*mobile_phone_number*’ (which is subsumed by the meaning of the query) but the contrary may not be correct. See Fig. 1 for an example of multilingual trees as input and of mappings as output of NuSM.

2.2. The semantic matching process

SMATCH matches its two input trees in a four-step process where the two first steps consist of formalising each input tree into a so-called *lightweight ontology* [1] while the two last steps perform the actual matching using background knowledge:

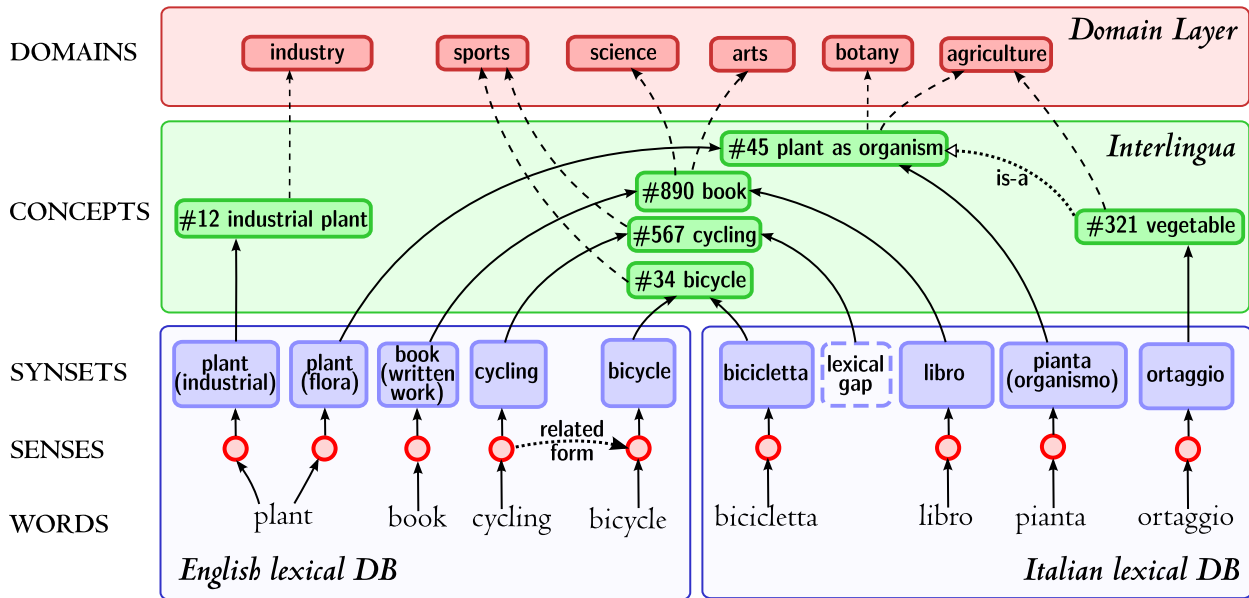


Fig. 3. Example English and Italian lexical databases with the interlingua layer providing language interoperability and the domain layer providing domain categories for concepts.

1. a lower layer of pluggable language-specific *lexical databases* that are WordNet-like lexical-semantic resources;
2. the *interlingua*: a language-independent but domain aware ontology of concepts where each concept is linked to its corresponding lexical entries in each language;
3. the *domain layer* that provides explicit domain information by associating interlingua concepts to domain categories.

The architecture of *lexical databases* is similar to that of Princeton WordNet [12], consisting of *lemmas* (i.e., dictionary forms of words of a language) associated to formally defined *word senses*. Synonymous senses are grouped together in synonym sets or *synsets*. Both senses and synsets are interconnected by lexical-semantic relations. Synsets represent an abstraction from the language-specific lexicon towards units of lexical meaning. Each lexical database is extensible by language-specific domain terms as described in Section 8.

The principal role of the language-independent *interlingua* layer is to serve as a bridge between language-specific lexical databases. Each synset in a lexical database is mapped to precisely one concept representing the corresponding language-independent unit of meaning. The opposite is not necessarily true. If a lexical database is incomplete or if a concept does not have a lexicalisation in a language, it will not be mapped to a synset in that language. This phenomenon is known as *lexical gap*. For example, there is no Italian word for *cycling*, which is paraphrased as ‘andare a bicicletta’, ‘to go on a bicycle’. In this case, the concept of *cycling* is connected to an English synset but not to any Italian synset as it has no lexicalisation in that language.

The interlingua acts as an interoperability layer across language-specific lexical entries, a feature that we use for cross-lingual matching. We consider two lexemes (words or expressions) in two different languages to be equivalent (to ‘mean the same’) if their synsets are connected to the same language-independent concept, such as the English synset of ‘book’ and the Italian synset of ‘libro’ in Fig. 3. Other relations (e.g., hypernymy, meronymy) can also be deduced across languages from concept relations (e.g., subsumption, part-of) of the interlingua.

The interlingua can also be used to provide abstraction from language-specific lexical meanings and to incorporate language-independent knowledge not provided by any underlying wordnet.

for example, through the introduction of new relations among concepts. This can be reused as additional background knowledge during the ontology matching task.

Finally, the *domain layer* annotates interlingua concepts with explicit domain information. Depending on the field of study (e.g., computational linguistics, knowledge representation, information retrieval), multiple formalisations exist for the notion of domain: lexical (as categories of words), semantic (as categories of concepts), or pragmatic (as categories of documents) [13]. We adopt the semantic approach and *define a domain as a labelled category of concepts*. Within the bounds of this general definition, different concrete models may be proposed and do exist (which we will shortly present); on this abstract level of definition, however, we do not impose any other constraints on how domains should be implemented as NuSM can be adapted to any of those models.

In our implementation of the background knowledge, as lexical databases we used wordnets from the *Open Multilingual WordNet* project.¹ Our interlingua layer is an extended and modified version of the Princeton WordNet synset graph. Our domain layer, finally, is a converted, language-independent version of the *Extended WordNet Domains* [14] resource. Section 8.1 gives more details on our implementation as well as a list of alternative, freely available components from which the three-layered background knowledge of NuSM can be built.

5. Language aware label parsing

The goal of the first step of semantic matching (cf. Section 2) is to formalise natural language ontology labels into propositional description logic formulas. This formal representation acts as a pivot and abstracts away various aspects of heterogeneity common in natural language such as ambiguity, variations in phrasing, and also multilingualism. In NuSM, label parsing is carried out in three substeps: language detection, computation of the structure of the formula, and computation of the atoms of the formula.

¹ <http://compling.hss.ntu.edu.sg/omw/>.

Please cite this article in press as: G. Bella, et al., Language and domain aware lightweight ontology matching, *Web Semantics: Science, Services and Agents on the World Wide Web* (2017), <http://dx.doi.org/10.1016/j.websem.2017.03.003>.

<i>Growing</i> POS = noun L = growing	<i>of</i> POS = closed-class OP = \sqcap	<i>plants</i> POS = noun L = plant	<i>and</i> POS = closed-class OP = \sqcup	<i>animals</i> POS = noun L = animal
---	--	--	---	--

Fig. 4. Result of tokenisation (boxes), part-of-speech tagging ('POS'), lemmatisation ('L'), and operator mapping ('OP') on the label 'Growing of plants and animals'.

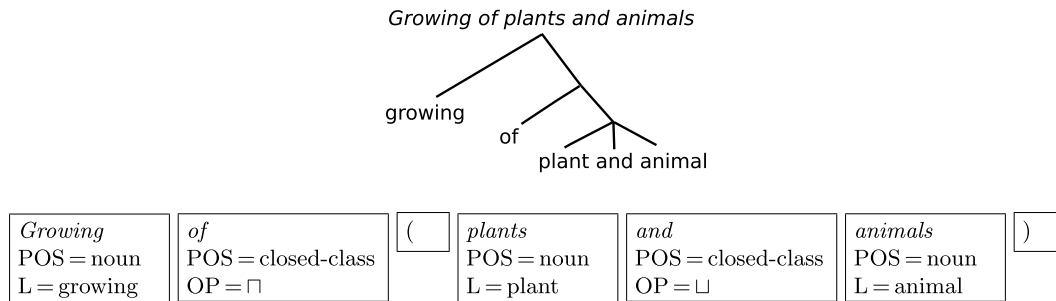


Fig. 5. Result of syntactic parsing and subsequent formula building on the label 'Growing of plants and animals': bracketing is introduced.

English	Italian	Spanish	Operator
except, non, without, ...	eccetto, escluso, non, senza, ...	excepto, sin, no, fuera de, salvo, ...	\neg
and, or, ' , ' , ...	e, ed, o, ' , ' , ...	y, e, o, u, ' , ' , ...	\sqcup
all other closed-class words			\sqcap
absence of closed-class word between two open-class words			\sqcap

Fig. 6. Mapping of closed-class words in labels to description logic operators (the list is incomplete). For example, the phrase 'English literature except poetry' will be mapped to $english \sqcap literature \sqcap \neg poetry$.

Furthermore, we account for phenomena such as paraphrasing and approximate translations that we observe to be more frequent across languages than in monolingual matching. As an illustration let us take the frequent case of English noun–noun compounds such as 'school transport' or 'river tourism', both taken from our EUROVOC evaluation corpus. In Latin languages (Italian, Spanish, French, etc.) these are typically translated as noun+adjective: 'trasporto scolastico', 'turismo fluviale'. In order to increase recall for such matches, we also retrieve meanings of *derivationally related words*, e.g., 'river–fluvial', 'school–scholastic–schooling–education', or 'cycling–bicycle' (cf. the example lexical gap on Fig. 3). This is why we are able to match 'river tourism' with 'turismo fluviale' and also 'cycling tourism' with 'turismo a bicicletta' in the example in Fig. 1.

6. Domain aware label disambiguation

The retrieval of all possible meanings for each atom results in highly ambiguous labels. Finding the most relevant concepts (meanings) for polysemous words is the role of the well-known NLP task of *word sense disambiguation* (WSD).

In semantic matching, sense disambiguation serves the purpose of avoiding erroneous mappings that decrease precision and may also decrease recall (if a wrong mapping is retained instead of the right one). More precisely, it eliminates word meanings that are incorrect in the given context: e.g., an equivalence relation between 'stock' as in a business context and 'broth' (or 'brodo' in Italian) in the context of cooking.

In NuSM a domain-based disambiguation method is used for two reasons: firstly, we observed that most real-world matching scenarios involve domain-specific ontologies. Secondly, our domain-based approach is independent of text length and is therefore well-suited to short ontology labels. The method relies on domain–concept mappings retrieved from the domain layer of our background knowledge (cf. Section 4) and is implemented based on

our previous results from [9]. We refer the reader to this paper for more details as well as for evidence on the efficiency of domain-based sense disambiguation on the block language of structured data.

An important feature of the method is that it is entirely language-independent since performed on the level of interlingua concepts. Costly language-specific WSD solutions are thus not necessary. The process is divided into two main steps:

1. *domain detection*: the relevant domain(s) of the ontology are made explicit through automated estimation;
2. *domain-based disambiguation*: the word meanings most relevant to the detected domain(s) are selected.

Just like for language detection, the domain detector assumes the input tree to be homogeneous with respect to its domain(s). This assumption is analogous to the *one-domain-per-discourse hypothesis* in computational linguistics that claims that 'multiple uses of a word in a coherent portion of text tend to share the same domain' [13, p. 28]. While the majority of real-world use cases does seem to fit our assumption, multidomain ontologies do exist, such as large national and international classifications in industry (NAICS³), commerce (SITC⁴), or libraries (UDC⁵). These, however, are invariably *faceted classifications* [16] that very clearly divide their contents domain-wise by top-level nodes. It is therefore straightforward to extract the domain-specific portions of such classifications as a preprocessing step and feed only homogeneous trees to the matcher.

³ The North American Industry Classification System. <https://www.census.gov/eos/www/naics/>.

⁴ Standard International Trade Classification. <https://unstats.un.org/unsd/trade/sitcrev4.htm>.

⁵ Universal Decimal Classification. <http://www.udcc.org>.

Lexical relations	Ontological relations	Mapped to
synonymy, similarity, derivationally related forms, pertainymy	identical concepts	\leftrightarrow
hypernymy, hyponymy, holonymy, meronymy	is-a, attribute-value, part-whole, substance, membership	\leftarrow or \rightarrow
antonymy		\neg

Fig. 7. Mapping of lexical relations (first column) and language-independent ontological relations (second column) to propositional logic operators used by the SAT-based matcher. Antonymy is used exclusively for monolingual matching. For hierarchical relations (such as *is-a* or *part-of*) transitivity is taken into account.

8.1. Reuse of existing multilingual knowledge bases

The three-layered lexical-semantic architecture we defined in Section 4 is approximated, to various extents, by several already existing and often freely available resources. The main differences among them reside in the details of their formal models and in the way they are populated (through manual effort, algorithmically, semi-automatically).

In earlier efforts such as EuroWordNet [18], MultiWordNet [19], or the Multilingual Central Repository [20], language-specific wordnets are constructed semi-automatically. Cross-lingual interoperability is provided by mapping non-English synsets to their English Princeton WordNet (PWN) counterparts. In other words, the English synset graph itself serves as the interlingua. This means that most of these multilingual resources inherit both the lexical coverage and the Anglo-Saxon lexical-semantic bias of PWN [21,22]. Nevertheless, they were and still are enormously popular and have often served as a basis for further efforts in building multilingual lexical-semantic resources.

A more recent and more extensive project is *BabelNet* [17]. It is automatically built from existing wordnets, including the *Open Multilingual WordNet* collection, as well as from other resources such as *Wikipedia*, *Wiktionary*, or *OmegaWiki*. Currently it covers 271 languages. Its interlingua layer, composed of *Babel synsets*, was enriched with a large number of synsets and relations, and is thus different in terms of organisation and richness from earlier resources.

Another recent project, developed at the University of Trento independently and in parallel with *BabelNet*, is the *Universal Knowledge Core* [16]. The UKC is a multi-layered knowledge resource with its lexical and interlingua layers corresponding to those described above. The current UKC supports 38 languages. Like *BabelNet*, it was mostly populated from existing freely available wordnets, principally from *Open Multilingual WordNet*. However, it also includes manually added linguistic knowledge, from single lexical entries to entire wordnets (e.g., the Mongolian wordnet [23]). The relative importance of manual extension and curation in the objective of maintaining a high-quality resource is a distinguishing feature of the UKC with respect to parallel solutions.

For our research we used the UKC as the underlying knowledge resource of NuSM, due to its off-line availability as well as practical reasons of pre-existing integration with our systems. However, as *BabelNet* and the UKC share the same architecture (as depicted in Fig. 3), the former could also be plugged in and used with NuSM, with the solutions explained in this paper remaining applicable and relevant. The modular architecture of NuSM makes the development of a *BabelNet* connector a relatively easy task and one to be investigated in the future.

The domain layer, as defined above, is also instantiated in several existing resources. The first such resource was *WordNet Domains* [24] where domains are labelled with strings and are defined as sets of Princeton WordNet synsets. *Extended WordNet Domains* [14] builds on the former but defines domains as fuzzy sets, i.e., domain-concept relations are annotated with weights. Finally, *WordNet* itself defines *domain term categories* that are represented not as labels but as synsets themselves but that, however, only categorise a small subset of the WordNet synsets.

For NuSM we implemented the domain layer as a modified version of *Extended WordNet Domains* where domains are mapped to concepts of the UKC interlingua as opposed to Princeton WordNet synsets, as described in detail in [9].

8.2. Extension by new wordnets

A large number of language-specific wordnets are downloadable from the web, in most cases under one of the common free licences.⁶ These resources, even though not always encoded in the same file formats, tend to follow the logical structure of PWN and so can be reused for our purposes. Often developed through research or community efforts, these resources offer variable levels of lexical coverage. If a wordnet does not exist at all for a language or its coverage is deemed inadequate, various automated [25] and manual (expert-sourced or crowd-sourced [26]) enrichment methods are applicable. While the presentation of these methods extends beyond the scope of this paper, we mention as an illustration that the *Open Multilingual Wordnet* project so far managed to integrate 34 wordnets, while its extended version, generated through automated methods from Wiktionary data, integrates 150 languages [25].

The mapping of the synsets of the new wordnet to interlingua concepts is straightforward to automate for all wordnets that are synset-aligned with PWN, which is generally the case. The interconnection of language-specific wordnets is still an actively researched topic and is one of the major tasks undertaken by the *Global WordNet Association*.⁷

In the case of the UKC, as explained in Section 8.1, the interlingua is a modified version of the PWN synset graph, and a mapping resource between the two graphs is constantly maintained. For a new language – say, Spanish – we perform the mapping as shown in the pseudocode below. Note that the same method can also be applied to multilingual lexical databases other than the UKC as long as they share the architecture described above.

Algorithm 1 Adding a new wordnet (sub)graph

```

1: procedure ADDNEWWORDNET(newWordNet)
2:   for newSynset in TRAVERSEBFS(newWordNet) do
3:     pwnSynset  $\leftarrow$  MAPToPWN(newSynset)
4:     if pwnSynset = None then
5:       concept  $\leftarrow$  ATTACHSYNSET(newSynset)
6:     else
7:       concept  $\leftarrow$  MAPToCONCEPT(pwnSynset)
8:     CONNECT(newSynset, ukcConcept)
9:   function ATTACHSYNSET(synset)
10:    newConcept  $\leftarrow$  CREATECONCEPT()
11:    for parent in GETHYPERNYMS(synset) do
12:      pwnParent  $\leftarrow$  MAPToPWN(parent)
13:      conceptParent  $\leftarrow$  MAPToCONCEPT(pwnParent)
14:      ADDISARELATION(newConcept, conceptParent)
15:      if domain  $\leftarrow$  GETDOMAIN(synset) = None then
16:        domain  $\leftarrow$  GETDOMAIN(conceptParent)
17:      MAPToDOMAIN(newConcept, domain)
18:   return newConcept

```

⁶ E.g., from the Global WordNet Association or from Open Multilingual Wordnet.

⁷ <http://globalwordnet.org>.

By *incompleteness of background knowledge* we refer to missing words, meanings, or semantic relations. Wordnets are domain-independent resources that lack specialised domain terminology. As already discussed in Section 4, the Spanish and Italian resources we used for our evaluations have about one-third or one-fourth the lexical coverage of Princeton WordNet. We quantified the effect of lexical incompleteness on matching scores in Section 10.3.

By *alternate wordings* we refer to the capacity of language to express the same (or very similar) meaning in several ways. It may seem trivial to point out that natural language phrases and expressions often cannot correctly be translated word-by-word from one language to another. In the context of multilingual classification labels, often written in a more controlled and semi-formal language due to their normative use, the hypothesis is not as trivial but still holds according to our observations. An example is the English label 'Manufacturing' translated into Italian as 'Attività manifatturiere', i.e., 'Manufacturing activities'. The label formula computation method in SMATCH and NuSM cannot deal with approximate matches resulting from the presence of additional concepts (here: the concept of activity). Google Translate, on the other hand, is built to be able to perform phrase-level statistical translation and thus tends to be more robust in such cases. In contrast, the rule-based Apertium lacks a phrase-level statistical translation capability, leading to significantly lower matching scores.

Finally, by *limitations of NLP* we cover a wide range of linguistic processing errors often due to the inherently difficult task of parsing short text labels. Mainstream NLP tools such as machine learning models for part-of-speech tagging, parsing, etc., tend to be trained on longer conventional text such as newswire or Wikipedia articles and, hence, tend to underperform on short ontology labels. While NLP in NuSM was tuned to short labels, we observed that among NLP-related mistakes by far the most frequently recurring ones were committed by the syntactic parser, resulting in incorrect bracketing in label formulas. These mistakes can partly be attributed to the weakness of our parsing logic, partly to the inherent ambiguity of short labels (e.g., the label 'floor and wall covering' could be correctly parsed both as $(\text{floor} \sqcup \text{wall}) \sqcap \text{covering}$ and as $\text{floor} \sqcup (\text{wall} \sqcap \text{covering})$, the former of which representing best the intended meaning).

We found the following to be the most typical mistakes made by Google Translate on our evaluation corpora.

- alternate wordings;
- committing on wrong word meanings;
- training anomalies;
- syntactic parsing mistakes;
- cumulative mistakes in non-English language pairs.

While to a different extent, *alternate wordings* are also a problem for statistical machine translation. While Google Translate is able to provide phrase-level translations which makes it inherently more robust to the phenomenon of alternate wordings than NuSM that operates on the word- and multiword-level, it nevertheless has its limitations and cannot translate, e.g., 'Psicopatologia' (meaning *psychopathology*) into 'Abnormal psychology', its English equivalent in the UDC corpus.

Committing on wrong word meanings is caused by the machine translator needing to commit on the meaning of a polysemous input word in order to produce a single piece of translated text as output. For example, the Italian label 'Lavoro e fatica', literally meaning 'Work and fatigue', is translated by Google into 'Work and effort', since 'effort' is indeed one of the meanings of 'fatica'. The two English words 'effort' and 'fatigue' having distinct meanings, translation-based matching fails. NuSM is not concerned by this problem as it does not try to disambiguate meanings before

matching happens: both meanings of 'fatica' are attempted to be matched.

Training anomalies are due to the fact that state-of-the-art statistical machine translation is to a large extent based on sentence- or word-aligned parallel corpora used as training material, often obtained through automated processing. Errors in the original content or in the preprocessing algorithms lead to strange translation mistakes such as '(psicotecnica)' (appearing within parentheses in the UDC corpus and meaning *psychotechnics*) being translated into '(Psycho)' (meaning *psychopath*), or 'politica dell'informazione' (meaning *information policy*) into 'political information'.

By *syntactic parsing mistakes* we refer to choosing the wrong parsing, especially when the phrase structure is ambiguous, e.g., 'Concetti e leggi generali' (meaning *General concepts and laws*) is translated into 'Concepts and general laws'. This results in an incorrect label formula being built by the matcher tool.

Cumulative mistakes in non-English language pairs refer to a phenomenon of an increased number of matching errors when none of the two input languages is English. When matching a Spanish tree against an Italian tree, both need to be translated into English, resulting in a higher probability of translation errors appearing compared to the case where one of the trees is in English.⁹ NuSM does not suffer from this effect and can even take advantage of the linguistic proximity of languages for improved results. This can be observed in our evaluations (those without OOV words, in order to eliminate the bias due to lower lexical coverage) where NuSM generally obtained better scores for the Spanish–Italian language pair (both being Romance languages sharing a similar morphology and syntax) while GoogleSM performed relatively worse.

9.2. Combination methods

Several approaches to combining matching techniques are possible; here we introduce two possibilities:

- as a simple fusion post-processor on mappings;
- different matchers on different inputs.

In simple fusion the two matchers are used as black boxes and only their output mappings are combined. This is the architecture depicted in Fig. 8. This solution is the simplest both conceptually and implementation-wise, but is also more limited in the kinds of combinations it supports.

Discrimination of matchers based on the label allows for a more efficient matching of deep single-domain or multidomain classifications. In such classifications, the deeper in the hierarchy labels are found the more domain-specific they tend to be. Since the most popular statistical machine translators offer good results on domain-agnostic or moderately domain-specific texts, a possible combination technique is to run translation-based matching on labels close to the root, and NuSM with its background knowledge extended by domain terminology applied on deeper levels.

While the second method appears to be a promising research direction, we have so far only experimented with the first one. The combined matcher architecture we considered is shown in Fig. 8. The output of both SMATCH and NuSM are mappings in the form of $(\text{source-node}, \text{target-node}, \text{rel})$ where rel is a semantic relation out of $\mathcal{R} = \{ \equiv, \sqsubset, \sqsupset, \emptyset \}$. The *combiner* component is a function

$$f_c : \mathcal{R} \times \mathcal{R} \rightarrow \mathcal{R}$$

⁹ Let us note that the problem would still remain if direct translation were applied from one language to another (e.g., from Spanish to Italian), as Google Translate always uses English as a pivot, resulting in two subsequent translations being performed.

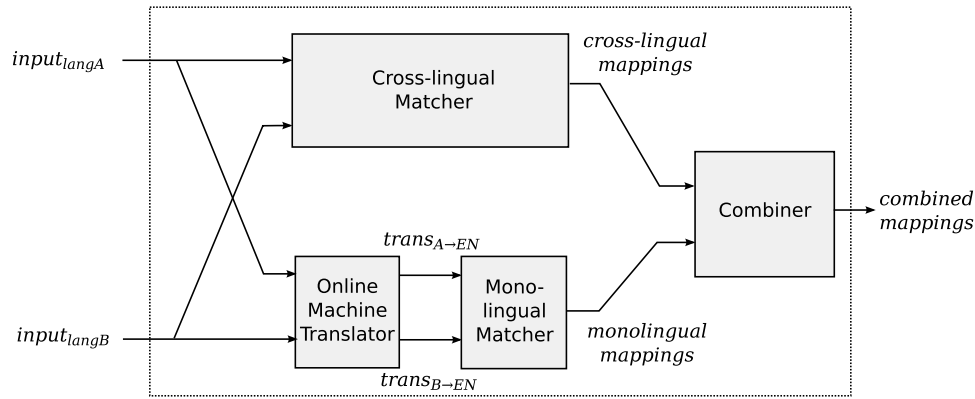


Fig. 8. A combined matcher that generates its output using the mappings output by its two component matchers.

f_V	\equiv	\sqsubset	\sqsupset	\emptyset
\equiv	\equiv	\equiv	\equiv	\equiv
\sqsubset	\equiv	\sqsubset	\emptyset	\sqsubset
\sqsupset	\equiv	\emptyset	\sqsupset	\sqsupset
\emptyset	\equiv	\sqsubset	\sqsupset	\emptyset

f_\wedge	\equiv	\sqsubset	\sqsupset	\emptyset
\equiv	\equiv	\sqsubset	\sqsupset	\emptyset
\sqsubset	\sqsubset	\sqsubset	\emptyset	\emptyset
\sqsupset	\sqsupset	\emptyset	\sqsupset	\emptyset
\emptyset	\emptyset	\emptyset	\emptyset	\emptyset

Fig. 9. Definition of combining functions f_{\vee} (left) and f_{\wedge} (right) based on the mapping relations output by NuSM and by the machine translation-based SMATCH. In case the two inputs do not match, f_{\vee} returns the stronger of the two while f_{\wedge} returns the weaker one. In case of equally strong but contradictory input mappings (e.g., \square and $\neg \square$) both functions behave in a conservative manner and output \emptyset (no relation).

taking for each source–target node pair the output relations of both matchers as input and returning a new relation rel_c :

$$rel_c := f_c(rel_{\text{NuSM}}, rel_{\text{MT}}).$$

We considered two possible f_c combining functions:

- a ‘greedy’ or ‘OR-like’ function f_{\vee} that emits the stronger of the two input mapping relations;
- a ‘conservative’ or ‘AND-like’ function f_{\wedge} that emits the weaker of the two input mapping relations.

The precise definitions of both functions are given in the tables in Fig. 9. Note that in the case of semantic matching neither SMATCH nor NuSM associates confidence scores to its mappings so the combiner can only rely on mapping relations.¹⁰

In a real-world use case the choice of combining function will be determined by practical needs: if priority is given to reducing the number of false positives then f_{\wedge} should be used as it is designed to be ‘conservative’ and return fewer false positives, thereby increasing precision at the cost of lower recall. On the other hand, f_{\vee} should be used in order to increase recall at the price of decreased precision, as it is defined so as to decrease the number of false negatives. As in real-world use cases, as also reflected by our evaluation results, precision tends to be much higher than recall, hence optimising by f_{\vee} generally leads to an improved F-measure and can thus be used as an all-purpose combiner.

10. Evaluation and discussion

The objective of our evaluations was not only to provide precision and recall figures on NuSM itself but mainly to compare the performance of our method to that of state-of-the-art

machine-translation-based monolingual matching. We did not include highly domain-specific ontologies in our evaluation corpus, nor did we extend NuSM with specialised domain terminology, in order to avoid an unfair comparison to general purpose machine translation tools that are not extensible in the same manner. In other words, to a certain extent, our evaluation results downplay one of the strengths of our approach, that is, the incremental adaptability of the background knowledge to the matching task.

We set up four separate evaluation scenarios: (1) comparison of NuSM to machine-translation-based matching over our full evaluation corpora; (2) the same comparison but over corpora guaranteed not to contain out-of-vocabulary words (i.e., not covered by our background knowledge), in order to get an idea of the effect of lexical incompleteness (or the lack thereof) on matching results; (3) evaluation of the fusion matcher; (4) evaluation of the relevance of word sense disambiguation for ontology matching.

10.1. Evaluation method

Our evaluations were performed on three language pairs: English–Spanish, English–Italian, and Spanish–Italian. As multilingual evaluation corpora we used the Universal Decimal Classification (UDC) and a randomly chosen subset of the EUROVOC vocabulary,¹¹ both available in several languages. Statistics on our evaluation corpora can be found in Fig. 10.

The two corpora we used are significantly different – and thus complement each other well – in terms of label size and complexity: EUROVOC labels tend to be very short (average length of 2.3 tokens) while UDC labels are longer (5.3 tokens). EUROVOC thus presents a use case where labels are syntactically very simple, yet their meanings are harder to identify due to reduced context. UDC, on the other hand, provides richer labels both syntactically and contextually, while the longer label presents a greater challenge for parsing with a proportionally higher probability of mistakes.

Earlier evaluations we had performed (such as those in our previous work [33]) showed that deep nesting of nodes in input classification hierarchies causes a significant drop in matching recall and thus introduces a bias into the results that makes any difference between label parsing methods appear as less pronounced. In order to focus on evaluating label matching and to avoid any interference in our results stemming from structure-level matching, we decided to ‘flatten’ the tree of the UDC corpus into a list of top-level nodes. EUROVOC is already a flat list of terms so it did not need to be transformed in any way. Thus, the – otherwise unchanged – step 2

¹⁰ For other kinds of matchers that only output equivalence mappings with confidence scores, f would need to compute a fusion score from the two input confidence scores.

¹¹ EUROVOC: the EU's multilingual thesaurus (eurovoc.europa.eu), UDC: Universal Decimal Classification (udcc.org). Of the latter we used Main Tables 1–7, excluding table 5 because it largely consists of Latin botanical and zoological named entities.

Fig. 10. Corpora used for evaluation.

Please cite this article in press as: G. Bella, et al., Language and domain aware lightweight ontology matching, *Web Semantics: Science, Services and Agents on the World Wide Web* (2017), <http://dx.doi.org/10.1016/j.websem.2017.03.003>.

Matcher	EUROVOC	UDC
NuSM	6,435	25,750
GoogleSM	6,220	23,054

(a) Total number of mappings created by each matcher on each corpus.

f_{\wedge}	EUROVOC	UDC
$\equiv \mapsto \sqsubset, \sqsupset$	133	404
$\equiv \mapsto \emptyset$	395	496
$\sqsubset, \sqsupset \mapsto \emptyset$	2,138	15,504

(b) Number of mappings modified by f_{\wedge} .

f_v	EUROVOC	UDC
$\sqsubset, \sqsupset \mapsto \equiv$	133	404
$\emptyset \mapsto \equiv$	395	496
$\emptyset \mapsto \sqsubset, \sqsupset$	2,125	15,431

(c) Number of mappings modified by f_v .

Fig. 13. Statistics on the influence of combining functions on mappings.

scenarios of low precision, such as very short and highly ambiguous labels.

11. Related work

A domain aware technique is used for ontology matching in [37]. Here, the matching configuration is highly asymmetric as the target ontology is DBpedia with a huge amount of nodes. The objective is thus to filter the contents of DBpedia to entities similar to the input. Furthermore, domains are not predefined with respect to lexical meanings but are computed with respect to the matching task in an unsupervised and ad-hoc manner from structural information taken from DBpedia, such as instance-class or superclass relations. The ontology labels themselves do not play any role in the computation of domains and the technique is entirely language-agnostic. In conclusion, this method is usable in combination with an existing matcher tool in scenarios where the target ontology is large and provides ample structural context for robust domain definitions. Our approach, in contrast, also works on small inputs in symmetric matching scenarios and assumes that the semantics of nodes are in a large part contained within the labels as opposed to structure and relations. This is the case of lightweight ontologies that are the main target of NuSM.

The dominant approach to cross-lingual ontology matching is to translate ontology labels to a common target language, thereby reducing the problem to monolingual (most often English-to-English) matching. State-of-the-art matchers rely either on Bing or the Google Translate API, including AML [4] and LogMap [6], the two tools that performed best in the *Multifarm* cross-lingual matching tasks of the 2014 Ontology Alignment Evaluation Initiative [36], the most authoritative evaluation effort for ontology matchers. Likewise, oft-cited publications on multilingual and cross-lingual matching [5,7,8] all propose methods that rely on some form of translation, using either online services or dictionaries either on the level of whole labels or on the level of individual words.

The idea of using interconnected lexical databases for cross-lingual matching also appears in the recent paper [34], in the context of a comparison between BabelNet and Google Translate as online word-level translation services. Beyond this initial similarity, our approach is conceptually different. The matching technique described in [34] uses BabelNet and Google Translate not as multilingual knowledge bases but, again, as online translator

services that retrieve all possible *translations of words* (i.e., lemmas) appearing in labels. A simple form of meaning-level reasoning is introduced by telling the online service to augment the set of returned lemmas by synonyms. Matching is thus performed on the word level of a chosen pivot language. We, on the other hand, perform matching directly on the language-independent level of concepts where beyond synonymy we are able to exploit a richer set of concept relations such as *subsumption* or *part-of*.

The ontology label matching problem can be reformulated as that of *textual similarity* or *entailment*: if $a \rightarrow b$ (textual entailment) then $a \sqsubseteq b$ (subsumption label mapping). The approach taken by NuSM can thus be regarded as a knowledge- and logic-based cross-lingual textual entailment operation (with the added complexity of formalising labels in the context of their ancestor nodes). Cross-lingual textual entailment has received some interest in the last years. The backbone of state-of-the-art solutions is usually a statistical approach (e.g., machine learning on parallel corpora [38], cross-lingual distributional semantics [39]) or – in most cases – a simple machine translation to a pivot language (English) [39,40]. They are sometimes combined with a knowledge-based approach for handling cases of monolingual synonymy and polysemy [40]. These methods are not fundamentally different from what we evaluated by combining Google Translate with English-to-English semantic matching, and present similar strengths and weaknesses: on the one hand, they can achieve good performance when the underlying machine translator, word embedding, or machine learning model is of high quality and is appropriate to the matching task. On the other hand, high quality is reached through access to large parallel or comparable training corpora, while appropriateness to the matching task requires the same corpora to be close enough to the domains to which the input belongs. Another particularity of the statistical approaches to textual similarity or entailment is their tendency to gloss over small differences that fundamentally change the meaning of a phrase. For example, the two phrases ‘*cereals and rice*’ and ‘*cereals except rice*’, the likes of which often appear in classifications, tend to be found very similar by statistical methods while they will be properly handled by NuSM that is able correctly to convert the two phrases into strictly non-matching formulas.

12. Conclusions

We have presented a new approach to semantic ontology matching that uses natively language and domain aware techniques, relying on off-line multilingual NLP and lexical-semantic resources. The results we obtained confirmed the viability of the method. When compared to the state-of-the-art cross-lingual matching technique based using two different machine translation tools, the three approaches turned out to score roughly similarly in terms of precision. In terms of recall, Google Translate reached equivalent to slightly better scores (+0%–15% depending on the language pair) while the Apertium machine translator fared much worse (–15%–20%).

We found our slightly lower scores with respect to Google to be partially due to the incompleteness of our local multilingual resources (both concerning lexical coverage and NLP processes). Indeed, with complete lexical coverage the differences in recall between the two methods are greatly reduced and our method takes the lead by 2%–6% on the Spanish–Italian language pair. This is a significant observation considering the fact that in the case of our matcher the issue of lexical incompleteness is under the control of the user: coverage issues can be – and are expected to be – addressed through the enrichment of background lexical-semantic knowledge by domain terminology and facts.

Finally, we built a fusion matcher that exploits the differences between the knowledge- and translation-based approaches by

Please cite this article in press as: G. Bella, et al., Language and domain aware lightweight ontology matching, Web Semantics: Science, Services and Agents on the World Wide Web (2017), <http://dx.doi.org/10.1016/j.websem.2017.03.003>.

- [29] Paul Buitelaar, Bogdan Sacaleanu, Extending synsets with medical terms, in: *Proceedings of the 1st International WordNet Conference*, January 21–25, o.A., Mysore, India, 2002.
- [30] S.E. Wright, G. Budin, *Handbook of terminology management, Application-oriented Terminology Management*, J. Benjamins, ISBN: 9789027221551, 2001. <https://books.google.it/books?id=UYm7XvBXm7QC>.
- [31] Bernardo Magnini, Manuela Speranza, Merging global and specialized linguistic ontologies, in: *Proceedings of the Workshop Ontolex-2002 Ontologies and Lexical Knowledge Bases, LREC-2002*, 2002, pp. 43–48.
- [32] Antonio Toral, Monica Monachini, Claudia Soria, Montse Cuadros, German Rigau, Wauter Bosma, Piek Vossen, Linking a domain thesaurus to WordNet and conversion to WordNet-LMF, in: *Proceedings of Second International Conference on Global Interoperability for Language Resources, ICGL2010*, Hong Kong, 2010.
- [33] Gábor Bella, Fausto Giunchiglia, Ahmed Ghassan Tawfik AbuRa'ed, Fiona McNeill, A multilingual ontology matcher, in: *Proceedings of OM-2015 located at ISWC 2015, CEUR-WS*, vol. 1545.
- [34] Mamoun Abu Helou, Matteo Palmonari, Mustafa Jarrar, Effectiveness of automatic translations for cross-lingual ontology mapping, *J. Artif. Intell. Res. (JAIR)* 55 (2016) 165–208.
- [35] Mikel L. Forcada, et al., Apertium: a free/open-source platform for rule-based machine translation, *Mach. Trans. (ISSN: 1573-0573)* 25 (2) (2011) 127–144. <http://dx.doi.org/10.1007/s10590-011-9090-0>.
- [36] Zlatan Dragisic, et al., Results of the ontology alignment evaluation initiative 2014, in: *ISWC 2014, Riva del Garda, Trentino, Italy, 2014*, pp. 61–104. http://ceur-ws.org/Vol-1317/oeai14_paper0.pdf.
- [37] Kristian Slabbekoorn, Laura Hollink, Geert-Jan Houben, Domain-Aware ontology matching, in: *International Semantic Web Conference, Springer, 2012*, pp. 542–558.
- [38] Yashar Mehdad, Matteo Negri, José Guilherme C. de Souza, FBK: cross-lingual textual entailment without translation, in: *Proceedings of the First Joint Conference on Lexical and Computational Semantics*, in: *SemEval '12, Association for Computational Linguistics, Stroudsburg, PA, USA, 2012*, pp. 701–705. <http://dl.acm.org/citation.cfm?id=2387636.2387755>.
- [39] Eneko Agirre, Carmen Banea, Daniel M. Cer, Mona T. Diab, Aitor Gonzalez-Agirre, Rada Mihalcea, German Rigau, Janyce Wiebe, Semeval-2016 task 1: semantic textual similarity, monolingual and cross-lingual evaluation, in: *Proceedings of the 10th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2016, San Diego, CA, USA, June 16–17, 2016*, pp. 497–511. <http://aclweb.org/anthology/S/S16/S16-1081.pdf>.
- [40] Julio Castillo, Marina Cardenas, Sagan: A machine translation approach for cross-lingual textual entailment, in: *Proceedings of the First Joint Conference on Lexical and Computational Semantics*, in: *SemEval '12, Association for Computational Linguistics, Stroudsburg, PA, USA, 2012*, pp. 721–726. <http://dl.acm.org/citation.cfm?id=2387636.2387759>.